

Nonasymptotic noisy lossy source coding

Victoria Kostina, Sergio Verdú

Dept. of Electrical Engineering, Princeton University, NJ 08544, USA

Abstract

This paper shows new general nonasymptotic achievability and converse bounds and performs their dispersion analysis for the lossy compression problem in which the compressor observes the source through a noisy channel. While this problem is asymptotically equivalent to a noiseless lossy source coding problem with a modified distortion function, nonasymptotically there is a noticeable gap in how fast their minimum achievable coding rates approach the rate-distortion function, providing yet another example in which at finite blocklengths one must put aside traditional asymptotic thinking.

Index Terms

Achievability, converse, finite blocklength regime, lossy data compression, noisy source coding, strong converse, dispersion, memoryless sources, Shannon theory.

I. INTRODUCTION

Consider a lossy compression setup in which the encoder has access only to a noise-corrupted version X of a source S , and we are interested in minimizing (in some stochastic sense) the distortion $d(S, Z)$ between the true source S and its rate-constrained representation Z (see Fig. 1). This problem arises if the object to be compressed is the result of an uncoded transmission over a noisy channel, or if it is observed data subject to errors inherent to the measurement system. Examples include speech in a noisy environment and photographs corrupted by noise introduced by the image sensor. Since we are concerned with reproducing the original noiseless

This work was supported in part by the National Science Foundation (NSF) under Grant CCF-1016625 and by the Center for Science of Information (CSol), an NSF Science and Technology Center, under Grant CCF-0939370.

This work was presented in part at the 2013 IEEE Information Theory Workshop [1].

source rather than preserving the noise, the distortion measure is defined with respect to the source.

The noisy source coding setting was introduced by Dobrushin and Tsybakov [2], who showed that when the goal is to minimize the average distortion, the noisy source coding problem is asymptotically equivalent to a surrogate noiseless source coding problem. Specifically, for a stationary memoryless source with single-letter distribution P_S observed through a stationary memoryless channel with single-input transition probability kernel $P_{X|S}$, the noisy rate-distortion function under a separable distortion measure is given by

$$R(d) = \min_{\substack{P_{Z|X}: \\ \mathbb{E}[d(S,Z)] \leq d \\ S-X-Z}} I(X; Z) \quad (1)$$

$$= \min_{\substack{P_{Z|X}: \\ \mathbb{E}[\bar{d}(X,Z)] \leq d}} I(X; Z) \quad (2)$$

where the surrogate per-letter distortion measure is

$$\bar{d}(a, b) = \mathbb{E}[d(S, b) | X = a]. \quad (3)$$

Therefore, in the limit of infinite blocklengths, the problem is equivalent to a conventional (noiseless) lossy source coding problem where the distortion measure is the conditional average of the original distortion measure given the noisy observation of the source. Berger [3, p.79] used the surrogate distortion measure (3) to streamline the proof of (2). Witsenhausen [4] explored the capability of distortion measures defined through conditional expectations such as in (3) to treat various so-called indirect rate distortion problems. Sakrison [5] showed that if both the source and its noise-corrupted version take values in a separable Hilbert space and the fidelity criterion is mean-squared error, then asymptotically, an optimal code can be constructed by first obtaining a minimum mean-square estimate of the source sequence based on its noisy observation, and then compressing the estimated sequence as if it were noise-free. Wolf and Ziv [6] showed that Sakrison's result holds even nonasymptotically, namely, that the minimum average distortion achievable in one-shot noisy compression of the object S can be written as

$$D^*(M) = \mathbb{E}[|S - \mathbb{E}[S|X]|^2] + \inf_{f,c} \mathbb{E}[|c(f(X)) - \mathbb{E}[S|X]|^2] \quad (4)$$

where the infimum is over all encoders $f: \mathcal{X} \mapsto \{1, \dots, M\}$ and all decoders $c: \{1, \dots, M\} \mapsto \widehat{\mathcal{M}}$, and \mathcal{X} and $\widehat{\mathcal{M}}$ are the alphabets of the channel output and the decoder output, respectively. It

is important to note that the choice of the mean-squared error distortion is crucial for the validity of the additive decomposition in (4). For vector quantization of a Gaussian signal corrupted by an additive independent Gaussian noise under weighted squared error distortion measure, Ayanoglu [7] found explicit expressions for the optimum quantization rule. Wolf and Ziv's result was extended to waveform vector quantization under weighted quadratic distortion measures and to autoregressive vector quantization under the Itakura-Saito distortion measure by Ephraim and Gray [8], and by Fisher, Gibson and Koo [9], who studied a model in which both encoder and decoder have access to the history of their past input and output blocks, allowing exploitation of inter-block dependence. Thus, the cascade of the optimal estimator followed by the optimal compressor achieves the minimum average distortion in those settings as well.

Under the logarithmic loss distortion measure [10], the noisy source coding problem reduces to the information bottleneck problem [11]. Indeed, in the information bottleneck method, the goal is to minimize $I(X; Z)$ subject to the constraint that $I(S; Y)$ exceeds a certain threshold. On the other hand, the noisy rate-distortion function under logarithmic loss is given by (2) in which $\mathbb{E}[\bar{d}(X, Z)]$ is replaced by $H(S|Y)$. The two optimization problems are equivalent. The solution to the information bottleneck problem thus acquires an operational meaning of the asymptotically minimum achievable noisy source coding rate under logarithmic loss.

In this paper, we give new nonasymptotic achievability and converse bounds for the noisy source coding problem, which generalize the noiseless source coding bounds in [12]. We observe that at finite blocklengths, the noisy coding problem is, in general, not equivalent to the noiseless coding problem with the surrogate distortion measure. Essentially, the reason is that taking the expectation in (3) dismisses the randomness introduced by the noisy channel, which nonasymptotically cannot be neglected. The additional randomness introduced by the channel slows down the rate of approach to the asymptotic fundamental limit in the noisy source coding problem compared to the asymptotically equivalent noiseless problem. Specifically, for noiseless source coding of stationary memoryless sources with separable distortion measure, we showed previously [12] (see also [13] for an alternative proof in the finite alphabet case) that the minimum number M of representation points compatible with a given probability ϵ of exceeding distortion threshold d can be written as

$$\log M^*(k, d, \epsilon) = kR(d) + \sqrt{k\mathcal{V}(d)}Q^{-1}(\epsilon) + o\left(\sqrt{k}\right) \quad (5)$$



Fig. 1. Noisy source coding.

where $\mathcal{V}(d)$ is the rate-dispersion function, explicitly identified in [12], and $Q^{-1}(\cdot)$ denotes the inverse of the complementary standard Gaussian cdf. In this paper, we show that for noisy source coding of a discrete stationary memoryless source over a discrete stationary memoryless channel under a separable distortion measure, $\mathcal{V}(d)$ in (5) is replaced by the noisy rate-dispersion function $\tilde{\mathcal{V}}(d)$, which can be expressed as

$$\tilde{\mathcal{V}}(d) = \mathcal{V}(d) + \lambda^{*2} \text{Var}(\mathbb{E}[d(S, Z^*)|S, X] - \mathbb{E}[d(S, Z^*)|X]) \quad (6)$$

where $\lambda^* = -R'(d)$, $\mathcal{V}(d)$ is the rate-dispersion function of the surrogate rate-distortion setup, and Z^* denotes the reproduction random variable that achieves the rate-distortion function (2). The difference between $\mathcal{V}(d)$ and $\tilde{\mathcal{V}}(d)$ is due to stochastic variability of the channel from S to X , which nonasymptotically cannot be neglected. Note, also, that $\tilde{\mathcal{V}}(d)$ cannot be expressed solely as a function of the source distribution and the surrogate distortion function.

The rest of the paper is organized as follows. After introducing the basic definitions in Section II, we proceed to show new general nonasymptotic converse and achievability bounds in Sections III and IV, respectively, along with their asymptotic analysis in Section V. Finally, the example of a binary source observed through an erasure channel is discussed in Section VI.

II. DEFINITIONS

Consider the one-shot setup in Fig. 1 where we are given the distribution P_S on the alphabet \mathcal{M} and the transition probability kernel $P_{X|S}: \mathcal{M} \rightarrow \mathcal{X}$. We are also given the distortion measure $d: \mathcal{M} \times \widehat{\mathcal{M}} \mapsto [0, +\infty]$, where $\widehat{\mathcal{M}}$ is the representation alphabet. An (M, d, ϵ) code is a pair of mappings $P_{U|X}: \mathcal{X} \mapsto \{1, \dots, M\}$ and $P_{Z|U}: \{1, \dots, M\} \mapsto \widehat{\mathcal{M}}$ such that $\mathbb{P}[d(S, Z) > d] \leq \epsilon$.

Define

$$\mathbb{R}_{S,X}(d) \triangleq \inf_{P_{Z|X}: \mathbb{E}[d(X,Z)] \leq d} I(X; Z) \quad (7)$$

where $\bar{d}: \mathcal{X} \times \widehat{\mathcal{M}} \mapsto [0, +\infty]$ is given by

$$\bar{d}(x, z) \triangleq \mathbb{E} [d(S, z) | X = x] \quad (8)$$

and, as in [12], assume that the infimum in (7) is achieved by some $P_{Z^*|X}$ such that the constraint is satisfied with equality. Noting that this assumption guarantees differentiability of $\mathbb{R}_{S,X}(d)$, denote

$$\lambda^* = -\mathbb{R}'_{S,X}(d) \quad (9)$$

Furthermore, define, for an arbitrary $P_{Z|X}$

$$\bar{d}_Z(s|x) \triangleq \mathbb{E} [d(S, Z) | X = x, S = s] \quad (10)$$

$$= \mathbb{E} [d(s, Z) | X = x] \quad (11)$$

where the expectation is with respect $P_{Z|XS} = P_{Z|X}$, and (11) follows from $S - X - Z$.

Definition 1 (noisy d-tilted information). *For $d > d_{\min}$, the noisy d-tilted information in $s \in \mathcal{M}$ given observation $x \in \mathcal{X}$ is defined as*

$$\tilde{j}_{S,X}(s, x, d) \triangleq D(P_{Z^*|X=x} \| P_{Z^*}) + \lambda^* \bar{d}_{Z^*}(s|x) - \lambda^* d \quad (12)$$

where $P_{Z^*|X}$ achieves the infimum in (7).

As we will see, the intuitive meaning of the noisy d-tilted information is the number of bits required to represent s within distortion d given observation x .

For the surrogate noiseless source coding problem, we know that almost surely ([14, Lemma 1.4], [15, Theorem 2.1])

$$j_X(x, d) = \iota_{X;Z^*}(x; Z^*) + \lambda^* \bar{d}(x, Z^*) - \lambda^* d \quad (13)$$

$$= D(P_{Z^*|X=x} \| P_{Z^*}) + \lambda^* \mathbb{E} [\bar{d}(x, Z^*) | X = x] - \lambda^* d \quad (14)$$

where $j_X(x, d)$ is the \bar{d} -tilted information in $x \in \mathcal{X}$ whose intuitive meaning is the number of bits required to represent x within distortion d and which is formally defined in [12, Definition 6]. From (12) and (14), we get

$$\tilde{j}_{S,X}(s, x, d) = j_X(x, d) + \lambda^* \bar{d}_{Z^*}(s|x) - \lambda^* \mathbb{E} [\bar{d}_{Z^*}(S|x) | X = x] \quad (15)$$

and

$$\mathbb{R}_{S,X}(d) = \mathbb{E} [\tilde{j}_{S,X}(S, X, d)] \quad (16)$$

$$= \mathbb{E} [j_X(X, d)] \quad (17)$$

If alphabets \mathcal{M} and \mathcal{X} are finite, for $(c, a) \in \mathcal{M} \times \mathcal{X}$, denote the partial derivatives

$$\dot{\mathbb{R}}_{S,X}(s, x, d) \triangleq \frac{\partial}{\partial P_{\bar{S}\bar{X}}(c, a)} \mathbb{R}_{\bar{S},\bar{X}}(d) \big|_{P_{\bar{S}\bar{X}}=P_{SX}} \quad (18)$$

The following theorem demonstrates that the noisy d-tilted information exhibits properties similar to those of the d-tilted information listed in [15, Theorem 2.2] (reproduced in Theorem 11 in Appendix A).

Theorem 1. *Assume that \mathcal{M} and \mathcal{X} are finite sets. Suppose that for all $P_{\bar{S}\bar{X}}$ in some Euclidean neighborhood of P_{SX} , $\text{supp}(P_{\bar{Z}^*}) = \text{supp}(P_{Z^*})$, where $\mathbb{R}_{\bar{S},\bar{X}}(d) = I(\bar{X}; \bar{Z}^*)$. Then*

$$\frac{\partial}{\partial P_{\bar{S}\bar{X}}(s, x)} \mathbb{E} [\tilde{j}_{\bar{S},\bar{X}}(S, X, d)] \big|_{P_{\bar{S}\bar{X}}=P_{SX}} = -\log e, \quad (19)$$

$$\dot{\mathbb{R}}_{S,X}(s, x, d) = \tilde{j}_{S,X}(s, x, d) - \log e, \quad (20)$$

$$\text{Var} \left(\dot{\mathbb{R}}_{S,X}(S, X, d) \right) = \text{Var} (\tilde{j}_{S,X}(S, X, d)). \quad (21)$$

Proof: Appendix B. ■

For a given distribution $P_{\bar{Z}}$ on $\widehat{\mathcal{M}}$ and $\lambda > 0$ define the transition probability kernel

$$dP_{\bar{Z}^*|X=x}(z) = \frac{dP_{\bar{Z}}(z) \exp(-\lambda \bar{d}(x, z))}{\mathbb{E} [\exp(-\lambda \bar{d}(x, \bar{Z}))]} \quad (22)$$

and define the function

$$\tilde{J}_{\bar{Z}}(s, x, \lambda) \triangleq D(P_{\bar{Z}^*|X=x} \| P_{\bar{Z}}) + \lambda \bar{d}_{\bar{Z}^*}(s|x) \quad (23)$$

$$= J_{\bar{Z}}(x, \lambda) + \lambda \bar{d}_{\bar{Z}^*}(s|x) - \lambda \mathbb{E} [\bar{d}_{\bar{Z}^*}(S|x) | X = x] \quad (24)$$

where

$$J_{\bar{Z}}(x, \lambda) \triangleq \log \frac{1}{\mathbb{E} [\exp(-\lambda \bar{d}(x, \bar{Z}))]} \quad (25)$$

where the expectation is with respect to the unconditional distribution of $P_{\bar{Z}}$. Similar to [15, (2.26)], we refer to the function

$$\tilde{J}_{\bar{Z}}(s, x, \lambda) - \lambda d \quad (26)$$

as the *generalized noisy d-tilted information*. As in the noiseless case, the generalized d-tilted information turns out to be relevant to the following optimization problem.

$$\mathbb{R}_{S,X;\bar{Z}}(d) \triangleq \inf_{\substack{P_{Z|X}: \\ \mathbb{E}[\bar{d}(X,Z)] \leq d}} D(P_{Z|X} \| P_{\bar{Z}} | P_X) \quad (27)$$

III. CONVERSE BOUNDS

Our main converse result for the nonasymptotic noisy lossy compression setting is the following lower bound on the excess distortion probability as a function of the code size.

Theorem 2 (Converse). *Any (M, d, ϵ) code must satisfy*

$$\epsilon \geq \inf_{P_{Z|X}} \sup_{\gamma \geq 0, P_{\bar{X}|\bar{Z}}} \left\{ \mathbb{P} \left[\iota_{\bar{X}|\bar{Z}\|X}(X; Z) + \sup_{\lambda \geq 0} \lambda(\bar{d}(S, Z) - d) \geq \log M + \gamma \right] - \exp(-\gamma) \right\} \quad (28)$$

where

$$\iota_{\bar{X}|\bar{Z}\|X}(x; z) \triangleq \log \frac{dP_{\bar{X}|\bar{Z}=z}}{dP_X}(x) \quad (29)$$

Proof: Let the encoder and decoder be the random transformations $P_{U|X}$ and $P_{Z|U}$, where U takes values in $\{1, \dots, M\}$. Denote

$$B_d(s) \triangleq \left\{ z \in \widehat{\mathcal{M}} : \bar{d}(s, z) \leq d \right\} \quad (30)$$

We have, for any $\gamma \geq 0$

$$\begin{aligned} & \mathbb{P} \left[i_{\bar{X}|\bar{Z}|X}(X; Z) + \sup_{\lambda \geq 0} \lambda(d(S, Z) - d) \geq \log M + \gamma \right] \\ &= \mathbb{P} \left[i_{\bar{X}|\bar{Z}|X}(X; Z) + \sup_{\lambda \geq 0} \lambda(d(S, Z) - d) \geq \log M + \gamma, d(S, Z) > d \right] \\ &+ \mathbb{P} \left[i_{\bar{X}|\bar{Z}|X}(X; Z) + \sup_{\lambda \geq 0} \lambda(d(S, Z) - d) \geq \log M + \gamma, d(S, Z) \leq d \right] \end{aligned} \quad (31)$$

$$= \mathbb{P}[d(S, Z) > d] + \mathbb{P}[i_{\bar{X}|\bar{Z}|X}(X; Z) \geq \log M + \gamma, d(S, Z) \leq d] \quad (32)$$

$$\leq \epsilon + \mathbb{P}[i_{\bar{X}|\bar{Z}|X}(X; Z) \geq \log M + \gamma] \quad (33)$$

$$\leq \epsilon + \frac{\exp(-\gamma)}{M} \mathbb{E} [\exp(i_{\bar{X}|\bar{Z}|X}(X; Z))] \quad (34)$$

$$\leq \epsilon + \frac{\exp(-\gamma)}{M} \sum_{u=1}^M \sum_{z \in \widehat{\mathcal{M}}} P_{Z|U}(z|u) \sum_{x \in \mathcal{X}} P_X(x) \exp(i_{\bar{X}|\bar{Z}|X}(x; z)) \quad (35)$$

$$= \epsilon + \frac{\exp(-\gamma)}{M} \sum_{u=1}^M \sum_{z \in \widehat{\mathcal{M}}} P_{Z|U}(z|u) \sum_{x \in \mathcal{X}} P_{\bar{X}|\bar{Z}}(x|z) \quad (36)$$

$$= \epsilon + \exp(-\gamma) \quad (37)$$

where

- (32) is by direct solution for the supremum;
- (34) is by Markov's inequality;
- (35) follows by upper-bounding

$$P_{U|X}(u|x) \leq 1 \quad (38)$$

for every $(x, u) \in \mathcal{M} \times \{1, \dots, M\}$.

Finally, (28) follows by choosing γ and $P_{\bar{X}|\bar{Z}}$ that give the tightest bound and $P_{Z|X}$ that gives the weakest bound in order to obtain a code-independent converse. ■

The following bound reduces the intractable (in high dimensional spaces) infimization over all conditional probability distributions $P_{Z|X}$ in Theorem 2 to infimization over the symbols of the output alphabet.

Corollary 1. Any (M, d, ϵ) code must satisfy

$$\epsilon \geq \sup_{\gamma \geq 0, P_{\bar{X}|\bar{Z}}} \left\{ \mathbb{E} \left[\inf_{z \in \widehat{\mathcal{M}}} \mathbb{P} \left[\iota_{\bar{X}|\bar{Z}\|X}(X; z) + \sup_{\lambda \geq 0} \lambda(d(S, z) - d) \geq \log M + \gamma | X \right] \right] - \exp(-\gamma) \right\} \quad (39)$$

Proof: We weaken (28) using

$$\begin{aligned} & \inf_{P_{Z|X}} \mathbb{P} \left[\iota_{\bar{X}|\bar{Z}\|X}(X; Z) + \sup_{\lambda \geq 0} \lambda(d(S, Z) - d) \geq \log M + \gamma \right] \\ &= \mathbb{E} \left[\inf_{P_{Z|X}} \mathbb{P} \left[\iota_{\bar{X}|\bar{Z}\|X}(X; Z) + \sup_{\lambda \geq 0} \lambda(d(S, Z) - d) \geq \log M + \gamma | X \right] \right] \end{aligned} \quad (40)$$

$$= \mathbb{E} \left[\inf_{z \in \widehat{\mathcal{M}}} \mathbb{P} \left[\iota_{\bar{X}|\bar{Z}\|X}(X; z) + \sup_{\lambda \geq 0} \lambda(d(S, z) - d) \geq \log M + \gamma | X \right] \right] \quad (41)$$

where we used $S - X - Z$. ■

In our asymptotic analysis, we will apply Corollary 1 with suboptimal choices of λ and $P_{\bar{X}|\bar{Z}}$.

Remark 1. If $\text{supp}(P_{Z^*}) = \widehat{\mathcal{M}}$ and $P_{X|S}$ is the identity mapping so that $d(S, z) = \bar{d}(X, z)$ almost surely, for every z , then Corollary 1 reduces to the noiseless converse in [12, Theorem 7] by using (13) after weakening (39) with $P_{\bar{X}|\bar{Z}} = P_{X|Z^*}$ and $\lambda = \lambda^*$.

IV. ACHIEVABILITY BOUNDS

Theorem 3 (Achievability). *There exists an (M, d, ϵ) code with*

$$\epsilon \leq \inf_{P_{\bar{Z}}} \int_0^1 \mathbb{E} [\mathbb{P}^M [\pi(X, \bar{Z}) > t | X]] dt \quad (42)$$

where $P_{X\bar{Z}} = P_X P_{\bar{Z}}$, and

$$\pi(x, z) = \mathbb{P} [d(S, z) > d | X = x] \quad (43)$$

Proof: The proof appeals to a random coding argument. Given M codewords (c_1, \dots, c_M) , the encoder f and decoder c achieving minimum excess distortion probability attainable with the given codebook operate as follows. Having observed $x \in \mathcal{X}$, the optimum encoder chooses

$$i^* \in \arg \min_i \pi(x, c_i) \quad (44)$$

with ties broken arbitrarily, so $f(x) = i^*$ and the decoder simply outputs c_{i^*} , so $c(f(x)) = c_{i^*}$.

The excess distortion probability achieved by the scheme is given by

$$\mathbb{P}[d(S, c(f(X))) > d] = \mathbb{E}[\pi(X, c(f(X)))] \quad (45)$$

$$= \int_0^1 \mathbb{P}[\pi(X, c(f(X))) > t] dt \quad (46)$$

$$= \int_0^1 \mathbb{E}[\mathbb{P}[\pi(X, c(f(X))) > t|X]] dt \quad (47)$$

Now, we notice that

$$1\{\pi(x, c(f(x))) > t\} = 1\left\{\min_{i \in 1, \dots, M} \pi(x, c_i) > t\right\} \quad (48)$$

$$= \prod_{i=1}^M 1\{\pi(x, c_i) > t\} \quad (49)$$

and we average (47) with respect to the codewords Z_1, \dots, Z_M drawn i.i.d. from $P_{\bar{Z}}$, independently of any other random variable, so that $P_{XZ_1 \dots Z_M} = P_X \times P_{\bar{Z}} \times \dots \times P_{\bar{Z}}$, to obtain

$$\int_0^1 \mathbb{E}\left[\prod_{i=1}^M \mathbb{P}[\pi(X, Z_M) > t|X]\right] dt = \int_0^1 \mathbb{E}[\mathbb{P}^M[\pi(X, \bar{Z}) > t|X]] dt \quad (50)$$

Since there must exist a codebook achieving excess distortion probability below or equal to the average over codebooks, (42) follows. ■

Remark 2. Notice that we have actually shown that the right-hand side of (42) gives the exact minimum excess distortion probability of random coding, averaged over codebooks drawn i.i.d. from P_Z .

Remark 3. In the noiseless case, $S = X$, so almost surely

$$\pi(X, z) = 1\{d(S, z) > d\} \quad (51)$$

and the bound in Theorem 3 reduces to the noiseless random coding bound in [12, Theorem 10].

The bound in (42) can be weakened to obtain the following result, which generalizes Shannon's bound for noiseless lossy compression (see e.g. [12, Theorem 1]).

Corollary 2. *There exists an (M, d, ϵ) code with*

$$\epsilon \leq \inf_{\gamma \geq 0, P_{Z|X}} \left\{ \mathbb{P}[d(S, Z) > d] + \mathbb{P}[\iota_{X;Z}(X; Z) > \log M - \gamma] + e^{-\exp(\gamma)} \right\} \quad (52)$$

where $P_{SXZ} = P_S P_{X|S} P_{Z|X}$.

Proof: Fix $\gamma \geq 0$ and transition probability kernel $P_{Z|X}$. Let $P_X \rightarrow P_{Z|X} \rightarrow P_Z$ (i.e. P_Z is the marginal of $P_X P_{Z|X}$), and let $P_{X\bar{Z}} = P_X P_Z$. We use the nonasymptotic covering lemma [16, Lemma 5] to establish

$$\mathbb{E} [\mathbb{P}^M [\pi(X, \bar{Z}) > t | X]] \leq \mathbb{P} [\pi(X, Z) > t] + \mathbb{P} [\iota_{X;Z}(X; Z) > \log M - \gamma] + e^{-\exp(\gamma)} \quad (53)$$

Applying (53) to (50) and noticing that

$$\int_0^1 \mathbb{P} [\pi(X, Z) > t] dt = \mathbb{E} [\pi(X, Z)] \quad (54)$$

$$= \mathbb{P} [d(S, Z) > d] \quad (55)$$

we obtain (52). ■

The following weakening of Theorem 3 is tighter than that in Corollary 2. It uses the generalized d-tilted information and is amenable to an accurate second-order analysis. See [15, Theorem 2.19] for a noiseless lossy compression counterpart.

Theorem 4 (Achievability, generalized d-tilted information). *Suppose that $P_{Z|X}$ is such that almost surely*

$$d(S, Z) = \bar{d}_Z(S|X) \quad (56)$$

Then there exists an (M, d, ϵ) code with

$$\begin{aligned} \epsilon \leq \inf_{\gamma, \beta, \delta, P_{\bar{Z}}} \Big\{ & \mathbb{E} \Big[\inf_{\lambda > 0} \Big\{ \mathbb{P} [D(P_{Z|X=x} \| P_{\bar{Z}}) + \lambda \bar{d}_Z(S|x) - \lambda(d - \delta) > \log \gamma - \log \beta | X] \\ & + \mathbb{P} [\bar{d}_Z(S|X) > d | X] \\ & + |1 - \beta \mathbb{P} [d - \delta \leq \bar{d}_Z(S|X) \leq d | X]|^+ \Big\} \Big] + e^{-\frac{M}{\gamma}} \Big\} \end{aligned} \quad (57)$$

Note that assumption (56) is satisfied, for example, by any constant composition code $P_{Z^k|X^k}$.

Proof: The bound in (42) implies that for an arbitrary $P_{\bar{Z}}$, there exists an (M, d, ϵ) code

with

$$\begin{aligned} \epsilon &\leq \int_0^1 \mathbb{E} [\mathbb{P}^M [\pi(X, \bar{Z}) > t|X]] dt \\ &\leq e^{-\frac{M}{\gamma}} \mathbb{E} \left[\min \left\{ 1, \gamma \int_0^1 \mathbb{P} [\pi(X, \bar{Z}) \leq t|X] dt \right\} \right] + \int_0^1 \mathbb{E} [|1 - \gamma \mathbb{P} [\pi(X, \bar{Z}) \leq t|X] dt|^+] \\ &\quad (58) \end{aligned}$$

$$\leq e^{-\frac{M}{\gamma}} + \int_0^1 \mathbb{E} [|1 - \gamma \mathbb{P} [\pi(X, \bar{Z}) \leq t|X] dt|^+] \quad (59)$$

where to obtain (58) we applied [17]

$$(1-p)^M \leq e^{-Mp} \leq e^{-\frac{M}{\gamma}} \min(1, \gamma p) + |1 - \gamma p|^+ \quad (60)$$

The first term in the right side of (58) is upper bounded using the following chain of inequalities.

$$\begin{aligned} &\int_0^1 |1 - \gamma \mathbb{P} [\pi(X, \bar{Z}) \leq t|X = x]|^+ \\ &\leq \int_0^1 |1 - \gamma \mathbb{E} [\exp(-\iota_{Z|X} \bar{Z}(x; Z)) \mathbf{1}_{\{\pi(x, Z) \leq t\}} |X = x]|^+ \end{aligned} \quad (61)$$

$$\leq \int_0^1 |1 - \gamma \mathbf{1}_{\{\pi(x) \leq t\}} \mathbb{E} [\exp(-\iota_{Z|X} \bar{Z}(x; Z))]|^+ \quad (62)$$

$$= \pi(x) + (1 - \pi(x)) |1 - \gamma \mathbb{E} [\exp(-\iota_{Z|X} \bar{Z}(x; Z))]|^+ \quad (63)$$

$$\leq \pi(x) + (1 - \pi(x)) |1 - \gamma \exp(-D(P_{Z|X=x} \| P_{\bar{Z}}))|^+ \quad (64)$$

$$\leq \pi(x) + |1 - \gamma \exp(-D(P_{Z|X=x} \| P_{\bar{Z}})) \mathbb{P} [\bar{\mathbf{d}}_Z(S|x) \leq d]|^+ \quad (65)$$

$$\leq \pi(x) + |1 - \gamma \exp(-D(P_{Z|X=x} \| P_{\bar{Z}})) \mathbb{P} [d - \delta \leq \bar{\mathbf{d}}_Z(S|x) \leq d]|^+ \quad (66)$$

$$\leq \pi(x) + |1 - \gamma \mathbb{E} [\exp(-g(S, x) - \lambda \delta)] \mathbb{P} [d - \delta \leq \bar{\mathbf{d}}_Z(S|x) \leq d]|^+ \quad (67)$$

$$\leq \pi(x) + \mathbb{P} [g(S, x) > \log \gamma - \log \beta - \lambda \delta] + |1 - \beta \mathbb{P} [d - \delta \leq \bar{\mathbf{d}}_Z(S|x) \leq d]|^+ \quad (68)$$

where

- in (62) we denoted

$$\pi(x) \triangleq \mathbb{P} [\bar{\mathbf{d}}_Z(S|x) > d] \quad (69)$$

where the probability is evaluated with respect to $P_{S|X=x}$, and observed using (56) that almost surely

$$\pi(X, Z) = \pi(X) \quad (70)$$

- (64) is by Jensen's inequality;
- in (64), we denoted

$$g(s, x) = D(P_{Z|X=x} \| P_{\bar{Z}}) + \lambda \bar{d}_Z(S|x) - \lambda d \quad (71)$$

- to obtain (68), we bounded

$$\gamma \exp(-g(S, x)) \geq \begin{cases} \beta & \text{if } g(S, x) \leq \log \gamma - \log \beta - \lambda \delta \\ 0 & \text{otherwise} \end{cases} \quad (72)$$

Taking the expectation of (68) and recalling (58), (57) follows. ■

V. ASYMPTOTIC ANALYSIS

In this section, we pass from the single shot setup of Sections III and IV to a block setting by letting the alphabets be Cartesian products $\mathcal{M} = \mathcal{S}^k$, $\mathcal{X} = \mathcal{A}^k$, $\widehat{\mathcal{M}} = \widehat{\mathcal{S}}^k$, and we study the second order asymptotics in k of $M^*(k, d, \epsilon)$, the minimum achievable number of representation points compatible with the excess distortion constraint $\mathbb{P} [d(S^k, Z^k) > d] \leq \epsilon$. We make the following assumptions.

- (i) $P_{S^k X^k} = P_S P_{X|S} \times \dots \times P_S P_{X|S}$ and

$$d(s^k, z^k) = \frac{1}{k} \sum_{i=1}^k d(s_i, z_i) \quad (73)$$

- (ii) The alphabets \mathcal{S} , \mathcal{A} , $\widehat{\mathcal{S}}$ are finite sets.

- (iii) The distortion level satisfies $d_{\min} < d < d_{\max}$, where

$$d_{\min} = \inf \{d: \mathbb{R}_{S, X}(d) < \infty\} \quad (74)$$

and $d_{\max} = \inf_{z \in \widehat{\mathcal{S}}} \mathbb{E} [d(X, z)]$, where the expectation is with respect to the unconditional distribution of X .

- (iv) The function $\mathbb{R}_{S, X; Z^*}(d)$ (defined in (27)) is twice continuously differentiable in some neighborhood of P_X .

The following result is obtained via a technical second order analysis of Corollary 1 (Appendix C) and Theorem 4 (Appendix D).

Theorem 5 (Gaussian approximation). *For $0 < \epsilon < 1$,*

$$\log M^*(k, d, \epsilon) = kR(d) + \sqrt{k\tilde{\mathcal{V}}(d)Q^{-1}(\epsilon)} + o(\sqrt{k}) \quad (75)$$

$$\tilde{\mathcal{V}}(d) = \text{Var}(\tilde{j}_{S;X}(S, X, d)) \quad (76)$$

Remark 4. The rate-dispersion function of the surrogate noiseless problem is given by (see [12, (83)])

$$\mathcal{V}(d) = \text{Var}(j_X(X, d)) \quad (77)$$

where $j_X(X, d)$ is defined in (13). To verify that the decomposition (6) indeed holds, which implies that $\tilde{\mathcal{V}}(d) > \mathcal{V}(d)$ unless there is no noise, write

$$\tilde{\mathcal{V}}(d) = \text{Var}(j_X(X, d) + \lambda^* \bar{d}_{Z^*}(S|X) - \lambda^* \mathbb{E}[\bar{d}_{Z^*}(S|X)|X]) \quad (78)$$

$$\begin{aligned} &= \text{Var}(j_X(X, d)) + \lambda^{*2} \text{Var}(\bar{d}_{Z^*}(S|X) - \mathbb{E}[\bar{d}_{Z^*}(S|X)|X]) \\ &\quad + 2\lambda^* \text{Cov}(j_X(X, d), \bar{d}_{Z^*}(S|X) - \mathbb{E}[\bar{d}_{Z^*}(S|X)|X]) \end{aligned} \quad (79)$$

where the covariance is zero:

$$\begin{aligned} &\mathbb{E}[(j_X(X, d) - R(d))(\bar{d}_{Z^*}(S|X) - \mathbb{E}[\bar{d}_{Z^*}(S|X)|X])] \\ &= \mathbb{E}[(j_X(X, d) - R(d))\mathbb{E}[\bar{d}_{Z^*}(S|X) - \mathbb{E}[\bar{d}_{Z^*}(S|X)|X]]] \end{aligned} \quad (80)$$

$$= 0 \quad (81)$$

Remark 5. Generalizing the observation made by Ingber and Kochman [13] in the noiseless lossy compression setting, we note using Theorem 1 that the noisy rate-dispersion function admits the following representation:

$$\tilde{\mathcal{V}}(d) = \text{Var}(\dot{\mathbb{R}}_{S;X}(S, X, d)) \quad (82)$$

VI. EXAMPLE: ERASED FAIR COIN FLIPS

A. Erased fair coin flips

Let a binary equiprobable source be observed by the encoder through a binary erasure channel with erasure rate δ . The goal is to minimize the bit error rate with respect to the source. For

$\frac{\delta}{2} \leq d \leq \frac{1}{2}$, the rate-distortion function is given by

$$R(d) = (1 - \delta) \left(\log 2 - h \left(\frac{d - \frac{\delta}{2}}{1 - \delta} \right) \right) \quad (83)$$

where $h(\cdot)$ is the binary entropy function, and (83) is obtained by solving the optimization in (2) which is achieved by achieved by $P_{Z^*}(0) = P_{Z^*}(1) = \frac{1}{2}$ and

$$P_{X|Z^*}(a|b) = \begin{cases} 1 - d - \frac{\delta}{2} & b = a \\ d - \frac{\delta}{2} & b \neq a \neq ? \\ \delta & a = ? \end{cases} \quad (84)$$

where $a \in \{0, 1, ?\}$ and $b \in \{0, 1\}$, so

$$\tilde{J}_{S,X}(S, X, b, d) = i_{X;Z^*}(X; b) + \lambda^* d(S, b) - \lambda^* d \quad (85)$$

$$= -\lambda^* d + \begin{cases} \log \frac{2}{1 + \exp(-\lambda^*)} & \text{w.p. } 1 - \delta \\ \lambda^* & \text{w.p. } \frac{\delta}{2} \\ 0 & \text{w.p. } \frac{\delta}{2} \end{cases} \quad (86)$$

The rate-dispersion function is given by the variance of (86):

$$\tilde{\mathcal{V}}(d) = \delta(1 - \delta) \log^2 \cosh \left(\frac{\lambda^*}{2 \log e} \right) + \frac{\delta}{4} \lambda^{*2} \quad (87)$$

$$\lambda^* = -R'(d) = \log \frac{1 - \frac{\delta}{2} - d}{d - \frac{\delta}{2}} \quad (88)$$

Bounds to the minimum achievable rate exploiting the symmetry of the erased coin flips setting were shown in [12].

B. Erased fair coin flips: surrogate rate-distortion problem

According to (8), the distortion measure of the surrogate rate-distortion problem is given by

$$\bar{d}(1, 1) = \bar{d}(0, 0) = 0 \quad (89)$$

$$\bar{d}(1, 0) = \bar{d}(0, 1) = 1 \quad (90)$$

$$\bar{d}(?, 1) = \bar{d}(?, 0) = \frac{1}{2} \quad (91)$$

The d-tilted information is given by taking the expectation of (86) with respect to S :

$$J_X(X, d) = -\lambda^* d + \begin{cases} \log \frac{2}{1+\exp(-\lambda^*)} & \text{w.p. } 1 - \delta \\ \frac{\lambda^*}{2} & \text{w.p. } \delta \end{cases} \quad (92)$$

Its variance is equal to

$$\mathcal{V}(d) = \delta(1 - \delta) \log^2 \cosh \left(\frac{\lambda^*}{2 \log e} \right) \quad (93)$$

$$= \tilde{\mathcal{V}}(d) - \frac{\delta}{4} \lambda^{*2} \quad (94)$$

For convenience, denote

$$\left\langle \begin{smallmatrix} k \\ j \end{smallmatrix} \right\rangle = \sum_{i=0}^j \binom{k}{i} \quad (95)$$

with the convention $\left\langle \begin{smallmatrix} k \\ j \end{smallmatrix} \right\rangle = 0$ if $j < 0$ and $\left\langle \begin{smallmatrix} k \\ j \end{smallmatrix} \right\rangle = \left\langle \begin{smallmatrix} k \\ k \end{smallmatrix} \right\rangle$ if $j > k$.

Achievability and converse bounds for the ternary source with binary representation alphabet and the distortion measure in (89)–(91) are obtained as follows.

Theorem 6 (Converse, surrogate BES). *Any (k, M, d, ϵ) code must satisfy*

$$\epsilon \geq \sum_{j=0}^{\lfloor 2kd \rfloor} \binom{k}{j} \delta^j (1 - \delta)^{k-j} \left[1 - M 2^{-(k-j)} \left\langle \begin{smallmatrix} k-j \\ \lfloor kd - \frac{1}{2}j \rfloor \end{smallmatrix} \right\rangle \right]^+ \quad (96)$$

Proof: Fix a (k, M, d, ϵ) code. While j erased bits contribute $\frac{1}{2} \frac{j}{k}$ to the total distortion regardless of the code, the probability that $k - j$ nonerased bits lie within Hamming distance ℓ of their representation can be upper bounded using [12, Theorem 15]:

$$\mathbb{P} \left[(k - j) \bar{d}(X^{k-j}, Z^{k-j}) \leq \ell \mid \text{no erasures in } X^{k-j} \right] \leq M 2^{-k+j} \left\langle \begin{smallmatrix} k-j \\ \ell \end{smallmatrix} \right\rangle \quad (97)$$

We have

$$\begin{aligned} & \mathbb{P} [\bar{d}(X^k, Z^k) \leq d] \\ &= \sum_{j=0}^{\lfloor 2kd \rfloor} \mathbb{P}[j \text{ erasures in } X^k] \mathbb{P} \left[(k - j) \bar{d}(X^{k-j}, Z^{k-j}) \leq kd - \frac{1}{2}j \mid \text{no erasures in } X^{k-j} \right] \end{aligned} \quad (98)$$

$$\leq \sum_{j=0}^{\lfloor 2kd \rfloor} \binom{k}{j} \delta^j (1 - \delta)^{k-j} \min \left\{ 1, M 2^{-(k-j)} \left\langle \begin{smallmatrix} k-j \\ \lfloor kd - \frac{1}{2}j \rfloor \end{smallmatrix} \right\rangle \right\} \quad (99)$$

■

Theorem 7 (Achievability, surrogate BES). *There exists a (k, M, d, ϵ) code such that*

$$\epsilon \leq \sum_{j=0}^k \binom{k}{j} \delta^j (1 - \delta)^{k-j} \left(1 - 2^{-(k-j)} \left\langle \begin{matrix} k-j \\ \lfloor kd - \frac{1}{2}j \rfloor \end{matrix} \right\rangle \right)^M \quad (100)$$

Proof: Consider the ensemble of codes with M codewords drawn i.i.d. from the equiprobable distribution on $\{0, 1\}^k$. Every erased symbol contributes $\frac{1}{2k}$ to the total distortion. The probability that the Hamming distance between the nonerased symbols and their representation exceeds ℓ , averaged over the code ensemble is found as in [12, Theorem 16]:

$$\mathbb{P}[(k-j)d(X^{k-j}, C(f(X^{k-j}))) > \ell \mid \text{no erasures in } X^{k-j}] = \left(1 - 2^{-(k-j)} \left\langle \begin{matrix} k-j \\ \ell \end{matrix} \right\rangle \right)^M \quad (101)$$

where $C(m)$, $m = 1, \dots, M$ are i.i.d on $\{0, 1\}^{k-j}$. Averaging over the erasure channel, we have

$$\begin{aligned} & \mathbb{P}[d(S^k, C(f(X^k))) > d] \\ &= \sum_{j=0}^k \mathbb{P}[j \text{ erasures in } X^k] \mathbb{P}\left[(k-j)d(S^{k-j}, C(f(X^{k-j}))) > kd - \frac{1}{2}j \mid \text{no erasures in } X^{k-j}\right] \end{aligned} \quad (102)$$

$$= \sum_{j=0}^k \binom{k}{j} \delta^j (1 - \delta)^{k-j} \left(1 - 2^{-(k-j)} \left\langle \begin{matrix} k-j \\ \lfloor kd - \frac{1}{2}j \rfloor \end{matrix} \right\rangle \right)^M \quad (103)$$

Since there must exist at least one code whose excess-distortion probability is no larger than the average over the ensemble, there exists a code satisfying (100). ■

The bounds in [12, Theorem 32], [12, Theorem 33] and the approximation in Theorem 5 (with the remainder term equal to 0 and $\frac{\log k}{2k}$), as well as the bounds in Theorems 6 and 7 for the surrogate rate-distortion problem together with their Gaussian approximation, are plotted in Fig. 2. Note that:

- The achievability and converse bounds are extremely tight, even at short blocklengths, as evidenced by Fig. 3 where we magnified the short blocklength region;
- The dispersion for the original noisy setup and its noiseless surrogate counterpart is small enough that the third-order term matters.
- Despite the fact that the asymptotically achievable rate in both problems is the same, there is a very noticeable gap between their nonasymptotically achievable rates in the displayed region of blocklengths. For example, at blocklength 1000, the penalty over the rate-distortion function is 9% for erased coin flips and only 4% for the surrogate source coding problem.

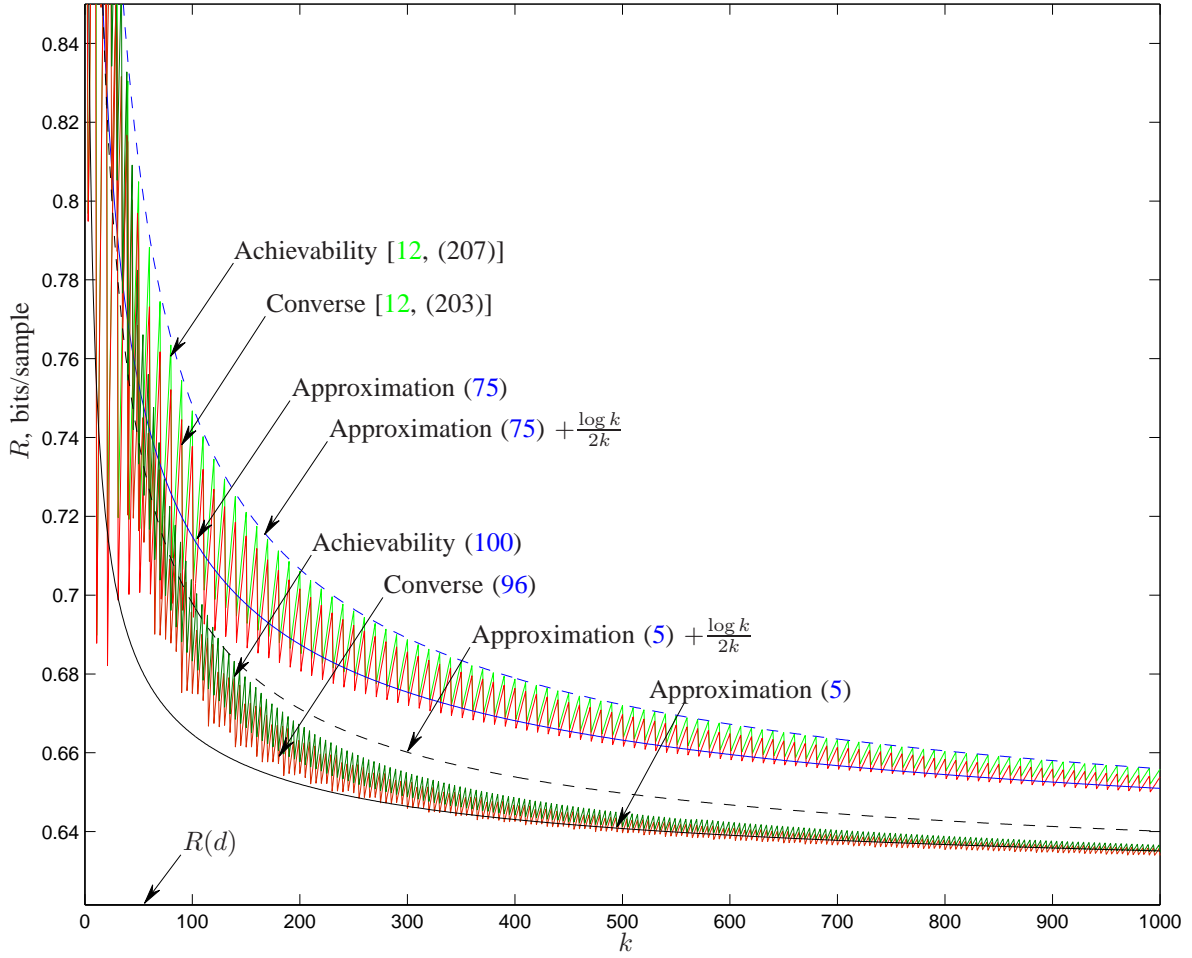


Fig. 2. Rate-blocklength tradeoff for the fair binary source observed through an erasure channel, as well as that for the surrogate problem, with $\delta = 0.1$, $d = 0.1$, $\epsilon = 0.1$.

APPENDIX A

AUXILIARY RESULTS

In this appendix, we state auxiliary results instrumental in the proof of Theorem 5.

Theorem 8 (Berry-Esseen CLT, e.g. [18, Ch. XVI.5 Theorem 2]). *Fix a positive integer k . Let W_i , $i = 1, \dots, k$ be independent. Then, for any real t*

$$\left| \mathbb{P} \left[\sum_{i=1}^k W_i > k \left(\mu_k + t \sqrt{\frac{V_k}{k}} \right) \right] - Q(t) \right| \leq \frac{B_k}{\sqrt{k}}, \quad (104)$$

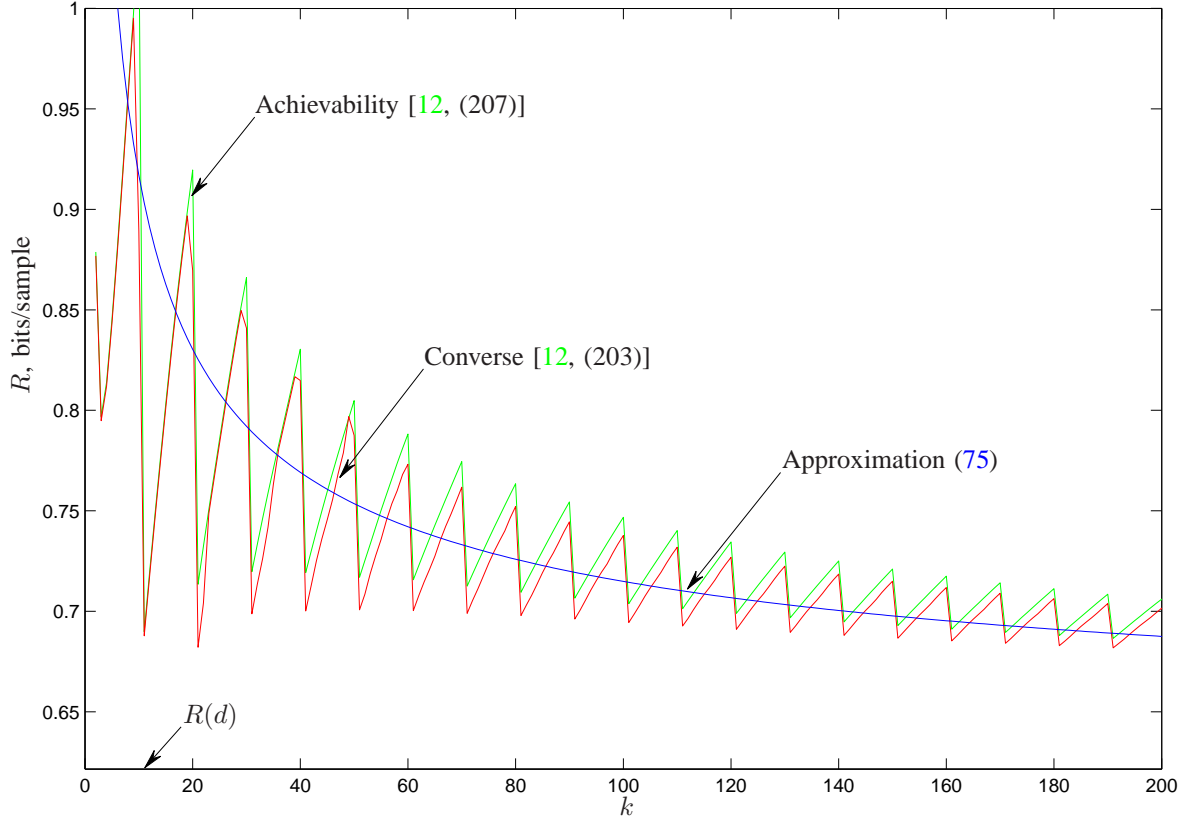


Fig. 3. Rate-blocklength tradeoff for the fair binary source observed through an erasure channel with $\delta = 0.1$, $d = 0.1$, $\epsilon = 0.1$ at short blocklengths.

where

$$\mu_k = \frac{1}{k} \sum_{i=1}^k \mathbb{E}[W_i] \quad (105)$$

$$V_k = \frac{1}{k} \sum_{i=1}^k \text{Var}(W_i) \quad (106)$$

$$T_k = \frac{1}{k} \sum_{i=1}^k \mathbb{E}[|W_i - \mu_i|^3] \quad (107)$$

$$B_k = \frac{c_0 T_k}{V_k^{3/2}} \quad (108)$$

and $0.4097 \leq c_0 \leq 0.5600$ ($0.4097 \leq c_0 < 0.4784$ for identically distributed W_i).

The second result deals with minimization of the cdf of a sum of independent random variables.

Let \mathcal{D} is a metric space with metric $d: \mathcal{D}^2 \mapsto \mathbb{R}^+$. Define the random variable Z on \mathcal{D} . Let $W_i, i = 1, \dots, k$ be independent conditioned on Z . Denote

$$\mu_k(z) = \frac{1}{k} \sum_{i=1}^k \mathbb{E} [W_i | Z = z] \quad (109)$$

$$V_k(z) = \frac{1}{k} \sum_{i=1}^k \text{Var} (W_i | Z = z) \quad (110)$$

$$T_k(z) = \frac{1}{k} \sum_{i=1}^k \mathbb{E} [|W_i - \mathbb{E} [W_i] |^3 | Z = z] \quad (111)$$

Let $\ell_1, \ell_2, L_1, F_1, F_2, V_{\min}$ and T_{\max} be positive constants. We assume that there exist $z^* \in \mathcal{D}$ and sequences μ_k^*, V_k^* such that for all $z \in \mathcal{D}$,

$$\mu_k^* - \mu_k(z) \geq \ell_1 d^2(z, z^*) - \frac{\ell_2}{k} \quad (112)$$

$$\mu_k^* - \mu_k(z^*) \leq \frac{L_1}{k} \quad (113)$$

$$|V_k(z) - V_k^*| \leq F_1 d(z, z^*) + \frac{F_2}{k} \quad (114)$$

$$V_{\min} \leq V_k(z) \quad (115)$$

$$T_k(z) \leq T_{\max} \quad (116)$$

Theorem 9 ([15, Theorem A.6.4]). *In the setup described above, under assumptions (112)–(116), for any $A > 0$, there exists a $K \geq 0$ such that, for all $|\Delta| \leq A 2\ell_1 T_{\max}^{\frac{1}{3}} V_{\min}^{\frac{5}{2}} F_1^{-2}$ and all sufficiently large k :*

$$\min_{z \in \mathcal{D}} \mathbb{P} \left[\sum_{i=1}^k W_i \leq k(\mu_k^* - \Delta) | Z = z \right] \geq Q \left(\Delta \sqrt{\frac{k}{V_k^*}} \right) - \frac{K}{\sqrt{k}} \quad (117)$$

The following two theorems summarize crucial properties of d-tilted information.

Theorem 10. [15, Theorem 2.1] *Fix $d > d_{\min}$. For P_Z^* -almost every z , it holds that*

$$j_X(x, d) = \imath_{X; Z^*}(x; z) + \lambda^* \bar{d}(x, z) - \lambda^* d \quad (118)$$

where $\lambda^* = -\mathbb{R}'_X(d)$, and $P_{XZ^*} = P_X P_{Z^*|X}$. Moreover,

$$\mathbb{R}_X(d) = \min_{P_{Z|X}} \mathbb{E} [\iota_{X;Z}(X; Z) + \lambda^* \bar{d}(X, Z)] - \lambda^* d \quad (119)$$

$$= \min_{P_{Z|X}} \mathbb{E} [\iota_{X;Z^*}(X; Z) + \lambda^* \bar{d}(X, Z)] - \lambda^* d \quad (120)$$

$$= \mathbb{E} [j_X(X, d)] \quad (121)$$

and for all $z \in \widehat{\mathcal{M}}$

$$\mathbb{E} [\exp \{ \lambda^* d - \lambda^* \bar{d}(X, z) + j_X(X, d) \}] \leq 1 \quad (122)$$

with equality for P_{Z^*} -almost every z .

For $a \in \mathcal{A}$, denote the partial derivatives

$$\dot{\mathbb{R}}_X(a) = \frac{\partial}{\partial P_{\bar{X}}(a)} \mathbb{R}_{\bar{X}}(d) \big|_{P_{\bar{X}}=P_X} \quad (123)$$

Theorem 11 ([15, Theorem 2.2]). Assume that \mathcal{X} is a finite set. Suppose that for all $P_{\bar{X}}$ in some neighborhood of P_X , $\text{supp}(P_{\bar{Z}^*}) = \text{supp}(P_{Z^*})$, where $\mathbb{R}_{\bar{X}}(d) = I(\bar{X}; \bar{Z}^*)$. Then,

$$\frac{\partial}{\partial P_{\bar{X}}(a)} \mathbb{E} [j_{\bar{X}}(X, d)] \big|_{P_{\bar{X}}=P_X} = \frac{\partial}{\partial P_{\bar{X}}(a)} \mathbb{E} [\iota_{\bar{X}}(X)] \big|_{P_{\bar{X}}=P_X} \quad (124)$$

$$= -\log e, \quad (125)$$

$$\dot{\mathbb{R}}_X(a) = j_X(a, d) - \log e, \quad (126)$$

$$\text{Var} \left(\dot{\mathbb{R}}_X(X) \right) = \text{Var} (j_X(X, d)). \quad (127)$$

Remark 6. The equality in (127) was first observed in [19].

Proof: Since by the assumption (118) particularized to $P_{\bar{X}}$ holds for P_{Z^*} -almost every z , we may write

$$\mathbb{E} [j_{\bar{X}}(X, d)] = \mathbb{E} [\iota_{\bar{X}; \bar{Z}^*}(X; Z^*)] - \mathbb{R}'_{\bar{X}}(d) \mathbb{E} [\bar{d}(X, Z^*) - d] \quad (128)$$

$$= \mathbb{E} [\iota_{\bar{X}; \bar{Z}^*}(X; Z^*)] \quad (129)$$

Therefore (in nats)

$$\frac{\partial}{\partial P_{\bar{X}}(a)} \mathbb{E} [j_{\bar{X}}(X, d)] \Big|_{P_{\bar{X}}=P_X} \quad (130)$$

$$= \frac{\partial}{\partial P_{\bar{X}}(a)} \mathbb{E} [\log P_{\bar{X}|\bar{Z}^*}(X; Z^*)] \Big|_{P_{\bar{X}}=P_X} - \frac{\partial}{\partial P_{\bar{X}}(a)} \mathbb{E} [\log P_{\bar{X}}(X)] \Big|_{P_{\bar{X}}=P_X} \quad (131)$$

$$= \frac{\partial}{\partial P_{\bar{X}}(a)} \mathbb{E} \left[\frac{P_{\bar{X}|\bar{Z}^*}(X; Z^*)}{P_{X|Z^*}(X; Z^*)} \right] \Big|_{P_{\bar{X}}=P_X} - \mathbb{E} \left[\frac{1}{P_X(X)} \frac{\partial}{\partial P_{\bar{X}}(a)} P_{\bar{X}}(X) \right] \Big|_{P_{\bar{X}}=P_X} \quad (132)$$

$$= \frac{\partial}{\partial P_{\bar{X}}(a)} 1 \Big|_{P_{\bar{X}}=P_X} - 1 \quad (133)$$

$$= -1 \quad (134)$$

This proves (125). To show (126), we invoke (125) to write

$$\dot{\mathbb{R}}_X(a) = \frac{\partial}{\partial P_{\bar{X}}(a)} \mathbb{E} [j_{\bar{X}}(\bar{X}, d)] \Big|_{P_{\bar{X}}=P_X} \quad (135)$$

$$= j_X(a, d) + \frac{\partial}{\partial P_{\bar{X}}(a)} \mathbb{E} [j_{\bar{X}}(X, d)] \Big|_{P_{\bar{X}}=P_X} \quad (136)$$

$$= j_X(a, d) - \log e \quad (137)$$

Finally, (127) is an immediate corollary to (126). ■

APPENDIX B

PROOF OF THEOREM 1

Denote for brevity $\bar{\lambda} = \lambda_{\bar{S}, \bar{X}}$. By the assumption (118) particularized to $P_X = P_{\bar{X}}$ holds for P_{Z^*} -almost every z , so using (15) we have almost surely

$$\begin{aligned} \tilde{j}_{\bar{S}, \bar{X}}(s, x, d) &= \iota_{\bar{X}; \bar{Z}^*}(x, Z^*) + \bar{\lambda} \mathbb{E} [\mathbf{d}(\bar{S}, Z^*) | X = x, Z^*] - \bar{\lambda} d \\ &\quad + \bar{\lambda} \bar{\mathbf{d}}_{\bar{Z}^*}(s|x) - \bar{\lambda} \mathbb{E} [\bar{\mathbf{d}}_{\bar{Z}^*}(\bar{S}|x) | \bar{X} = x] \end{aligned} \quad (138)$$

Therefore

$$\begin{aligned} \mathbb{E} [\tilde{j}_{\bar{S}, \bar{X}}(S, X, d)] &= \mathbb{E} [\iota_{\bar{X}; \bar{Z}^*}(X, Z^*)] + \bar{\lambda} \mathbb{E} [\bar{\mathbf{d}}_{Z^*}(\bar{S}|X)] - \bar{\lambda} d \\ &\quad + \bar{\lambda} \mathbb{E} [\bar{\mathbf{d}}_{\bar{Z}^*}(S|X)] - \bar{\lambda} \mathbb{E} [\bar{\mathbf{d}}_{\bar{Z}^*}(\bar{S}|X)] \end{aligned} \quad (139)$$

We make note of the following. In nats,

$$\begin{aligned} & \frac{\partial}{\partial P_{\bar{S}\bar{X}}(c, a)} \mathbb{E} [\iota_{\bar{X}; \bar{Z}^*}(X, Z^*)] \Big|_{P_{\bar{S}\bar{X}}=P_{SX}} \\ &= \frac{\partial}{\partial P_{\bar{S}\bar{X}}(c, a)} \mathbb{E} [\log P_{\bar{X}|\bar{Z}^*}(S; Z^*)] \Big|_{P_{\bar{S}\bar{X}}=P_{SX}} - \frac{\partial}{\partial P_{\bar{S}\bar{X}}(c, a)} \mathbb{E} [\log P_{\bar{X}}(X)] \Big|_{P_{\bar{S}\bar{X}}=P_{SX}} \end{aligned} \quad (140)$$

$$= \frac{\partial}{\partial P_{\bar{S}\bar{X}}(c, a)} \mathbb{E} \left[\frac{P_{\bar{X}|\bar{Z}^*}(X; Z^*)}{P_{X|Z^*}(X; Z^*)} \right] \Big|_{P_{\bar{S}\bar{X}}=P_{SX}} - \mathbb{E} \left[\frac{1}{P_X(X)} \frac{\partial}{\partial P_{\bar{S}\bar{X}}(c, a)} P_{\bar{X}}(X) \right] \Big|_{P_{\bar{S}\bar{X}}=P_{SX}} \quad (141)$$

$$= \frac{\partial}{\partial P_{\bar{S}\bar{X}}(c, a)} 1 \Big|_{P_{\bar{S}\bar{X}}=P_{SX}} - 1 \quad (142)$$

$$= -1 \quad (143)$$

Moreover,

$$\begin{aligned} & \frac{\partial}{\partial P_{\bar{S}\bar{X}}(c, a)} \mathbb{E} [\bar{d}_{Z^*}(\bar{S}|X)] \Big|_{P_{\bar{S}\bar{X}}=P_{SX}} \\ &= \sum_{x,z} P_{Z^*|X}(z|x) P_X(x) \sum_s d(s, z) \frac{\partial}{\partial P_{\bar{S}\bar{X}}(c, a)} P_{\bar{S}|\bar{X}}(s|x) \Big|_{P_{\bar{S}\bar{X}}=P_{SX}} \end{aligned} \quad (144)$$

$$= \bar{d}_{Z^*}(c|a) - \mathbb{E} [\bar{d}_{Z^*}(S|a)] \quad (145)$$

where we used

$$\frac{\partial}{\partial P_{\bar{S}\bar{X}}(c, a)} P_{\bar{S}|\bar{X}}(s|x) \Big|_{P_{\bar{S}\bar{X}}=P_{SX}} = \frac{1 \{x = a\}}{P_X(a)} (1 \{s = c\} - P_{S|X}(s|a)) \quad (146)$$

Similarly,

$$\begin{aligned} & \frac{\partial}{\partial P_{\bar{S}\bar{X}}(c, a)} \mathbb{E} [\bar{d}_{\bar{Z}^*}(S|X) - \bar{d}_{\bar{Z}^*}(\bar{S}|X)] \Big|_{P_{\bar{S}\bar{X}}=P_{SX}} \\ &= \frac{\partial}{\partial P_{\bar{S}\bar{X}}(c, a)} \sum_{x,z} P_{\bar{Z}^*|X}(z|x) P_X(x) \sum_s d(s, z) (P_{S|X}(s|x) - P_{\bar{S}|\bar{X}}(s|x)) \Big|_{P_{\bar{S}\bar{X}}=P_{SX}} \end{aligned} \quad (147)$$

$$= - \sum_{x,z} P_{Z^*|X}(z|x) P_X(x) \sum_s d(s, z) \frac{\partial}{\partial P_{\bar{S}\bar{X}}(c, a)} P_{\bar{S}|\bar{X}}(s|x) \Big|_{P_{\bar{S}\bar{X}}=P_{SX}} \quad (148)$$

$$= \mathbb{E} [\bar{d}_{Z^*}(S|a)] - \bar{d}_{Z^*}(c|a) \quad (149)$$

Assembling (143), (145), (149), it follows that

$$\frac{\partial}{\partial P_{\bar{S}\bar{X}}(c, a)} \mathbb{E} [\tilde{j}_{\bar{S}, \bar{X}}(S, X, d)] \Big|_{P_{\bar{S}\bar{X}}=P_{SX}} = -\log e \quad (150)$$

This proves (125). To show (126), we invoke (125) to write

$$\dot{\mathbb{R}}_{S,X}(a) = \frac{\partial}{\partial P_{\bar{S}\bar{X}}(c,a)} \mathbb{E} [\tilde{j}_{\bar{S},\bar{X}}(\bar{S}, \bar{X}, d)] \Big|_{P_{\bar{S}\bar{X}}=P_{SX}} \quad (151)$$

$$= j_{S,X}(c, a, d) + \frac{\partial}{\partial P_{\bar{S}\bar{X}}(c,a)} \mathbb{E} [j_{\bar{S},\bar{X}}(S, X, d)] \Big|_{P_{\bar{S}\bar{X}}=P_{SX}} \quad (152)$$

$$= j_{S,X}(a, d) - \log e \quad (153)$$

Finally, (127) is an immediate corollary to (126).

APPENDIX C

PROOF OF THE CONVERSE PART OF THEOREM 5

The proof consists of an asymptotic analysis of the bound in Theorem 2 with a careful choice of tunable parameters.

The following auxiliary result will be instrumental.

Lemma 1. *Let X_1, \dots, X_k be independent on \mathcal{A} and distributed according to P_X . For all k , it holds that*

$$\mathbb{P} \left[|\text{type}(X^k) - P_X|^2 > \frac{\log k}{k} \right] \leq \frac{|\mathcal{A}|}{\sqrt{k}} \quad (154)$$

Proof: By Hoeffding's inequality, similar to Yu and Speed [20, (2.10)]. ■

Let $P_{Z|X}: \mathcal{A} \mapsto \hat{\mathcal{S}}$ be a stochastic matrix whose entries are multiples of $\frac{1}{k}$. We say that the conditional type of z^k given x^k is equal to $P_{Z|X}$, $\text{type}(z^k|x^k) = P_{Z|X}$, if the number of a 's in x^k that are mapped to b in z^k is equal to the number of a 's in x^k times $P_{Z|X}(b|a)$, for all $(a, b) \in \mathcal{A} \times \hat{\mathcal{S}}$.

Let

$$\log M = kR(d) + \sqrt{k\tilde{\mathcal{V}}(d)Q^{-1}(\epsilon_k)} - \frac{1}{2} \log k - \log |\mathcal{P}_{[k]}| \quad (155)$$

where $\epsilon_k = \epsilon + o(1)$ will be specified in the sequel, and $\mathcal{P}_{[k]}$ denotes the set of all conditional k -types $\hat{\mathcal{S}} \rightarrow \mathcal{A}$.

We weaken the bound in (28) by choosing

$$P_{\bar{X}^k|\bar{Z}^k=z^k}(x^k) = \frac{1}{|\mathcal{P}_{[k]}|} \sum_{P_{X|Z} \in \mathcal{P}_{[k]}} \prod_{i=1}^k P_{X|Z=z_i}(x_i) \quad (156)$$

$$\lambda = k\lambda(x^k) = k\mathbb{R}'_{\text{type}(x^k)}(d) \quad (157)$$

$$\gamma = \frac{1}{2} \log k \quad (158)$$

By virtue of Theorem 2, the excess distortion probability of all (M, d, ϵ) codes where M is that in (155) must satisfy

$$\epsilon \geq \mathbb{E} \left[\min_{z^k \in \hat{\mathcal{S}}^k} \mathbb{P} \left[\iota_{\bar{X}^k|\bar{Z}^k\|X^k}(X^k; z^k) + k\lambda(X^k)(d(S^k, z^k) - d) \geq \log M + \gamma|X^k| \right] \right] - \exp(-\gamma) \quad (159)$$

We identify the typical set of channel outputs:

$$\mathcal{T}_k = \left\{ x^k \in \mathcal{A}^k : |\text{type}(x^k) - P_X|^2 \leq \frac{\log k}{k} \right\} \quad (160)$$

where $|\cdot|$ is the Euclidean norm.

We proceed to evaluate the minimum in in (159) for $x^k \in \mathcal{T}_k$.

For a given pair (x^k, z^k) , abbreviate

$$\text{type}(x^k) = P_{\bar{X}} \quad (161)$$

$$\text{type}(z^k|x^k) = P_{\bar{Z}|\bar{X}} \quad (162)$$

$$\lambda(x^k) = \lambda_{\bar{X}} \quad (163)$$

We define $P_{\bar{X}|\bar{Z}}$ through $P_{\bar{X}}P_{\bar{Z}|\bar{X}}$ and lower-bound the sum in (156) by the term containing $P_{\bar{X}|\bar{Z}}$, concluding that

$$\begin{aligned} \iota_{\bar{X}^k|\bar{Z}^k\|X^k}(x^k; z^k) + \lambda(x^k)(d(S^k, z^k) - d) &\geq kI(\bar{X}; \bar{Z}) + kD(\bar{X}\|X) \\ &\quad + \lambda_{\bar{X}} \left(\sum_{i=1}^k \bar{d}_{\bar{Z}}(S_i|x_i) - kd \right) - \log |\mathcal{P}_{[k]}| \end{aligned} \quad (164)$$

$$= \sum_{i=1}^k W_i + kD(\bar{X}\|X) - \log |\mathcal{P}_{[k]}| \quad (165)$$

where

$$W_i = I(\bar{X}; \bar{Z}) + \lambda_{\bar{X}} (\bar{d}_{\bar{Z}}(S_i|x_i) - d) \quad (166)$$

where $P_{\bar{S}\bar{X}\bar{Z}}(s, a, b) = P_{\bar{X}}(a)P_{\bar{S}|\bar{X}}(s|a)P_{\bar{Z}|\bar{X}}(b|a)$.

Conditioned on $X^k = x^k$, the random variables W_i are independent with (in the notation of Theorem 9 where $z = P_{\bar{Z}|\bar{X}}$)

$$\mu_k(P_{\bar{Z}|\bar{X}}) = I(\bar{X}; \bar{Z}) + \lambda_{\bar{X}} (\mathbb{E} [\bar{d}_{\bar{Z}}(S|\bar{X})] - d) \quad (167)$$

$$V_k(P_{\bar{Z}|\bar{X}}) = \lambda_{\bar{X}}^2 \mathbb{E} [\text{Var} (\bar{d}_{\bar{Z}}(S|\bar{X})|\bar{X})] \quad (168)$$

$$T_k(P_{\bar{Z}|\bar{X}}) = \lambda_{\bar{X}}^3 \mathbb{E} \left[\left| \bar{d}_{\bar{Z}}(S|\bar{X}) - \mathbb{E} [\bar{d}_{\bar{Z}}(S|\bar{X})|\bar{X}] \right|^3 \right] \quad (169)$$

Denote the backward conditional distribution that achieves $\mathbb{R}_{\bar{X}}(d)$ by $P_{\bar{X}|\bar{Z}^*}$. Write

$$\mu_k(P_{\bar{Z}|\bar{X}}) = I(\bar{X}; \bar{Z}) + \lambda_{\bar{X}} (\mathbb{E} [\bar{d}(\bar{X}, \bar{Z})] - d) \quad (170)$$

$$= \mathbb{E} [\iota_{\bar{X}; \bar{Z}^*}(\bar{X}; \bar{Z}) + \lambda_{\bar{X}} \bar{d}(\bar{X}, \bar{Z})] - \lambda_{\bar{X}} d + D(P_{\bar{X}|\bar{Z}} \| P_{\bar{X}|\bar{Z}^*} | P_{\bar{Z}}) \quad (171)$$

$$\geq \mathbb{R}_{\bar{X}}(d) + D(P_{\bar{X}|\bar{Z}} \| P_{\bar{X}|\bar{Z}^*} | P_{\bar{Z}}) \quad (172)$$

$$\geq \mathbb{R}_{\bar{X}}(d) + \frac{1}{2} |P_{\bar{X}|\bar{Z}} P_{\bar{Z}} - P_{\bar{X}|\bar{Z}^*} P_{\bar{Z}}|^2 \log e \quad (173)$$

where (172) is by Theorem 10, and (173) is by Pinsker's inequality. Similar to the proof of [15, (C.50)], we conclude that the conditions of Theorem 9 are satisfied.

Denote

$$a_k = \log M + \frac{1}{2} \log k + \log |\mathcal{P}_{[k]}| - kD(\bar{X} \| \mathbf{X}) - k\mathbb{R}_{\bar{X}}(d) \quad (174)$$

$$b_k = \log M + \frac{1}{2} \log k + \log |\mathcal{P}_{[k]}| - kD(\bar{X} \| \mathbf{X}) - k\mathbb{R}_{\mathbf{X}}(d) - c \log k \quad (175)$$

$$W_i^* = j_{\mathbf{X}}(X_i, d) - \mathbb{R}_{\mathbf{X}}(d) + \lambda_{\mathbf{X}} \bar{d}_{\mathbf{Z}^*}(S_i | X_i) - \lambda_{\mathbf{X}} \mathbb{E} [\bar{d}_{\mathbf{Z}^*}(S_i | X_i) | X_i] \quad (176)$$

where M is that in (155), and constant $c > 0$ will be identified later. Weakening (159) further,

we have

$$\epsilon \geq \mathbb{E} \left[\min_{P_{\bar{Z}|\bar{X}}} \mathbb{P} \left[\sum_{i=1}^k W_i \geq k\mathbb{R}_{\bar{X}}(d) + a_k |\text{type}(X^k) = P_{\bar{X}} \right] 1 \{X^k \in \mathcal{T}_k\} \right] - \frac{1}{\sqrt{k}} \quad (177)$$

$$\geq \mathbb{E} \left[\mathbb{P} \left[\lambda_{\bar{X}} \left(\sum_{i=1}^k \bar{d}_{\bar{Z}^*}(S_i|X_i) - kd \right) \geq a_k |\text{type}(X^k) = P_{\bar{X}} \right] 1 \{X^k \in \mathcal{T}_k\} \right] - \frac{K+1}{\sqrt{k}} \quad (178)$$

$$\geq \mathbb{E} \left[\mathbb{P} \left[\lambda_{\bar{X}} \left(\sum_{i=1}^k \bar{d}_{\bar{Z}^*}(S_i|X_i) - k\mathbb{E} [\bar{d}_{\bar{Z}^*}(S|\bar{X})] \right) \geq a_k |\text{type}(X^k) = P_{\bar{X}} \right] 1 \{X^k \in \mathcal{T}_k\} \right] - \frac{K_1 \log k + K + 1}{\sqrt{k}} \quad (179)$$

$$\geq \mathbb{E} \left[\mathbb{P} \left[\sum_{i=1}^k W_i^* \geq b_k |\text{type}(X^k) = P_{\bar{X}} \right] 1 \{X^k \in \mathcal{T}_k\} \right] - \frac{K_1 \log k + K + 2B + 1}{\sqrt{k}} \quad (180)$$

$$\geq \mathbb{P} \left[\sum_{i=1}^k W_i^* \geq b_k \right] - \mathbb{P} [X^k \notin \mathcal{T}_k] - \frac{K_1 \log k + K + 2B + 1}{\sqrt{k}} \quad (181)$$

$$\geq \mathbb{P} \left[\sum_{i=1}^k W_i^* \geq b_k \right] - \frac{K_1 \log k + K + 2B + |\mathcal{A}| + 1}{\sqrt{k}} \quad (182)$$

$$\geq \epsilon_k - \frac{K_1 \log k + K + 2B + B^* + |\mathcal{A}| + 1}{\sqrt{k}} \quad (183)$$

where

- (178) is by Theorem 9, and $K > 0$ is defined therein.
- To show (179), which holds for some $K_1 > 0$, observe that since

$$\mathbb{E} [\bar{d}_{\bar{Z}^*}(S|\bar{X})] = \mathbb{E} [\bar{d}(\bar{X}, \bar{Z}^*)] \quad (184)$$

$$= d \quad (185)$$

conditioned on x^k , both random variables $\lambda_{\bar{X}} \left(\sum_{i=1}^k \bar{d}_{\bar{Z}^*}(S_i|x_i) - kd \right)$

and $\lambda_{\bar{X}} \left(\sum_{i=1}^k \bar{d}_{\bar{Z}^*}(S_i|x_i) - k\mathbb{E} [\bar{d}_{\bar{Z}^*}(S|\bar{X})] \right)$ are zero mean. By the Berry-Esséen theorem

and the assumption that all alphabets are finite, there exists $B > 0$ such that

$$\begin{aligned} & \mathbb{P} \left[\lambda_{\bar{X}} \left(\sum_{i=1}^k \bar{d}_{\bar{Z}^*}(S_i|X_i) - kd \right) \geq a_k |\text{type}(X^k) = P_{\bar{X}} \right] \\ & \geq Q \left(\frac{a_k}{\lambda_{\bar{X}} \sqrt{k \text{Var}(\bar{d}_{\bar{Z}^*}(S|\bar{X})|\bar{X})}} \right) - \frac{B}{\sqrt{k}} \end{aligned} \quad (186)$$

$$\geq Q \left(\frac{a_k}{\lambda_{\bar{X}} \sqrt{k \text{Var}(\bar{d}_{\bar{Z}^*}(S|\bar{X})|\bar{X})}} \left(1 + a \sqrt{\frac{\log k}{k}} \right) \right) - \frac{B}{\sqrt{k}} \quad (187)$$

$$\geq Q \left(\frac{a_k}{\lambda_{\bar{X}} \sqrt{k \text{Var}(\bar{d}_{\bar{Z}^*}(S|\bar{X})|\bar{X})}} \right) - \frac{B}{\sqrt{k}} - K_1 \frac{\log k}{\sqrt{k}} \quad (188)$$

$$\begin{aligned} & \geq \mathbb{P} \left[\lambda_{\bar{X}} \left(\sum_{i=1}^k \bar{d}_{\bar{Z}^*}(S_i|X_i) - k \mathbb{E}[\bar{d}_{\bar{Z}^*}(S|\bar{X})] \right) \geq a_k |\text{type}(X^k) = P_{\bar{X}} \right] \\ & - \frac{2B}{\sqrt{k}} - K_1 \frac{\log k}{\sqrt{k}} \end{aligned} \quad (189)$$

where (187) for some scalar a is obtained by applying a Taylor series expansion to $\frac{1}{\sqrt{\text{Var}(\bar{d}_{\bar{Z}^*}(S|\bar{X})|\bar{X})}}$ in the neighborhood of typical $P_{\bar{X}}$, i.e. those types corresponding to $x^k \in \mathcal{T}_k$, and (188) invokes (see e.g. [15, (A.32)])

$$Q(x + \xi) \geq Q(x) - \frac{|\xi|^+}{\sqrt{2\pi}} \quad (190)$$

with $\xi \sim \frac{\log k}{\sqrt{k}}$ because $a_k = O(\sqrt{k \log k})$ for typical $P_{\bar{X}}$.

- (180) holds because due to Taylor's theorem, there exists $c > 0$ such that

$$\mathbb{R}_{\bar{X}}(d) \geq \mathbb{R}_X(d) + \sum_{a \in \mathcal{A}} (P_{\bar{X}}(a) - P_X(a)) \dot{\mathbb{R}}_X(a) - c |P_{\bar{X}} - P_X|^2 \quad (191)$$

$$= \mathbb{R}_X(d) + \frac{1}{k} \sum_{i=1}^k \dot{\mathbb{R}}_X(X_i) - \mathbb{E}[\dot{\mathbb{R}}_X(X)] - c |P_{\bar{X}} - P_X|^2 \quad (192)$$

$$= \frac{1}{k} \sum_{i=1}^k \mathcal{J}_X(X_i, d) - c |P_{\bar{X}} - P_X|^2 \quad (193)$$

$$\geq \frac{1}{k} \sum_{i=1}^k \mathcal{J}_X(X_i, d) - c \log k \quad (194)$$

where (193) uses (126), and (194) is by the definition (160) of the typical set of x^k 's.

- (182) is by Lemma 1.

- (183) applies the Berry-Esséen theorem to the sequence of i.i.d. random variables W_i^* whose Berry-Esséen ratio is denoted by B^* .

The result now follows by letting

$$\epsilon_k = \epsilon + \frac{K + 2B + B^* + |\mathcal{A}| + 1 + K_1 \log k}{\sqrt{k}} \quad (195)$$

in (155).

APPENDIX D

PROOF OF THE ACHIEVABILITY PART OF THEOREM 5

The proof consists of an asymptotic analysis of the bound in Theorem 4 with a careful choice of auxiliary parameters so that only the first term in (57) survives.

Let $P_{\bar{Z}^k} = P_{Z^{k*}} = P_{Z^*} \times \dots \times P_{Z^*}$, where Z^* achieves $\mathbb{R}_X(d)$, and let $P_{\bar{X}} = P_{\bar{X}^k} = P_{\bar{X}} \times \dots \times P_{\bar{X}}$, where $P_{\bar{X}}$ is the measure on \mathcal{X} generated by the empirical distribution of $x^k \in \mathcal{X}^k$:

$$P_{\bar{X}}(a) = \frac{1}{k} \sum_{i=1}^k 1\{x_i = a\} \quad (196)$$

We let $\mathcal{T} = \mathcal{T}_k$ in (160) so that by Lemma 1

$$P_{X^k}(\mathcal{T}_k^c) \leq \frac{|\mathcal{A}|}{\sqrt{k}} \quad (197)$$

so we will concern ourselves only with typical x^k .

Let $P_{Z^*|\bar{X}}$ be the transition probability kernel that achieves $\mathbb{R}_{\bar{X}, Z^*}(d)$, and let $P_{Z^{*k}|X^k} = P_{Z^*|X} \times \dots \times P_{Z^*|X}$. Let $P_{Z^k|X^k}$ be uniform on the conditional type which is closest to (in terms of Euclidean distance)

$$P_{Z|X=x}(z) = \frac{P_{Z^*}(z) \exp(-\lambda \bar{d}(x, z))}{\mathbb{E}[\exp(-\lambda \bar{d}(x, Z^*))]} \quad (198)$$

(cf. (22)) where

$$\lambda = \lambda_{\bar{X}} = -\mathbb{R}'_{\bar{X}, Z^*}(d - \xi) \quad (199)$$

$$\xi = \sqrt{\frac{a \log k}{k}} \quad (200)$$

for some $0 < a < 1$, so that (56) holds, and

$$\mathbb{E}[\bar{d}(\bar{X}, Z)] = d - \xi \quad (201)$$

where $P_{\bar{X}} \rightarrow P_{Z|X} \rightarrow P_Z$.

It follows by the Berry-Esséen Theorem that

$$\mathbb{P} \left[\sum_{i=1}^k \bar{d}_Z(S_i|X_i = x_i) > kd | X^k = x^k \right] \leq \frac{1}{\sqrt{2\pi a k^a \log k}} + \frac{B}{\sqrt{k}} \quad (202)$$

where where B is the maximum (over $x^k \in \mathcal{T}_k$) of the Berry-Esséen ratios for $\bar{d}_Z(S_i|X_i = x_i)$, and we used

$$Q(x) < \frac{1}{\sqrt{2\pi x}} e^{-\frac{x^2}{2}} \quad (203)$$

Again by the Berry-Esséen theorem, we have

$$\begin{aligned} & \mathbb{P} \left[kd - \tau \leq \sum_{i=1}^k \bar{d}_Z(S_i|X_i = x_i) \leq kd | X^k = x^k \right] \\ & \geq Q \left(\sqrt{k}\xi - \frac{\tau}{\sqrt{k}} \right) - Q \left(\sqrt{k}\xi \right) - \frac{2B}{\sqrt{k}} \end{aligned} \quad (204)$$

$$\geq \frac{\tau}{\sqrt{2\pi k}} e^{-\frac{k\xi^2}{2}} - \frac{2B}{\sqrt{k}} \quad (205)$$

$$= \frac{1}{\sqrt{k}} \left(\frac{\tau}{\sqrt{2\pi k^a}} - 2B \right) \quad (206)$$

$$\geq \frac{b}{\sqrt{k}} \quad (207)$$

where (206) holds because for $x, y \geq 0$

$$Q(x+y) \geq Q(x) - \frac{y}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} \quad (208)$$

and (207) follows by choosing

$$\tau = (2B + b)\sqrt{2\pi k^a} \quad (209)$$

for some $b > 0$, k large enough and some $\tau > 0$.

We now proceed to evaluate the first term in (57).

$$D(P_{Z^k|X^k=x^k} \| P_{Z^*} \times \dots \times P_{Z^*}) \leq kD(Z \| Z^*) + kH(Z) - kH(Z|\bar{X}) + |\mathcal{A}||\hat{\mathcal{S}}| \log(k+1) \quad (210)$$

$$= kD(P_{Z|X} \| P_{Z^*} | P_{\bar{X}}) + |\mathcal{A}||\hat{\mathcal{S}}| \log(k+1) \quad (211)$$

where to obtain (220) we used the type counting lemma [21, Lemma 2.6]. Therefore

$$g(s^k, x^k) \triangleq D(P_{Z^k|X^k=x^k} \| P_{Z^*} \times \dots \times P_{Z^*}) + k\lambda_{\bar{X}} \bar{d}_{Z^k}(s^k|x^k) - k\lambda_{\bar{X}} d \quad (212)$$

$$= kD(P_{Z|X} \| P_{Z^*} | P_{\bar{X}}) + \lambda_{\bar{X}} \sum_{i=1}^k \mathbb{E} [\bar{d}_Z(s_i|\bar{X})] - k\lambda_{\bar{X}} d + |\mathcal{A}| |\hat{\mathcal{S}}| \log(k+1) \quad (213)$$

$$= \mathbb{E} [J_{Z^*}(\bar{X}, \lambda_{\bar{X}})] - k\lambda_{\bar{X}} d + k\lambda_{\bar{X}} \sum_{i=1}^k \mathbb{E} [\bar{d}_Z(s_i|\bar{X})] - \lambda_{\bar{X}} \mathbb{E} [\bar{d}_Z(S|\bar{X})] + |\mathcal{A}| |\hat{\mathcal{S}}| \log(k+1) \quad (214)$$

$$\leq k\mathbb{E} [J_{Z^*}(\bar{X}, \lambda_{\bar{X}, Z^*}^*)] - k\lambda_{\bar{X}, Z^*}^* d + \lambda_{\bar{X}} \sum_{i=1}^k \mathbb{E} [\bar{d}_Z(s_i|\bar{X})] - \lambda_{\bar{X}} \mathbb{E} [\bar{d}_Z(S|\bar{X})] + |\mathcal{A}| |\hat{\mathcal{S}}| \log(k+1) + L \log k \quad (215)$$

where to show (215) recall that by the assumption $\mathbb{R}_{\bar{X}, Z^*}(d)$ is twice continuously differentiable, so there exists $a > 0$ such that

$$\lambda - \lambda_{\bar{X}, Z^*}^* = \mathbb{R}'_{\bar{X}, Z^*}(d - \xi) - \mathbb{R}'_{\bar{X}, Z^*}(d) \quad (216)$$

$$\leq a\xi \quad (217)$$

Since $\lambda_{\bar{X}, Z^*}^*$ is a maximizer of $\mathbb{E} [J_{Z^*}(\bar{X}, \lambda)] - \lambda d$ (see e.g. [12, (261)])

$$\frac{\partial}{\partial \lambda} \mathbb{E} [J_{Z^*}(\bar{X}, \lambda)] |_{\lambda=\lambda_{\bar{X}, Z^*}^*} = d \quad (218)$$

the first term in the Taylor series expansion of $\mathbb{E} [J_{Z^*}(\bar{X}, \lambda)] - \lambda d$ in the vicinity of $\lambda_{\bar{X}, Z^*}^*$ vanishes, and we conclude that there exists L such that

$$\mathbb{E} [J_{Z^*}(\bar{X}, \lambda)] - \lambda d \geq \mathbb{E} [J_{Z^*}(\bar{X}, \lambda_{\bar{X}, Z^*}^*)] - \lambda_{\bar{X}, Z^*}^* d - L\xi^2 \quad (219)$$

Moreover, according to [12, Lemma 5], there exist $C_2, K_2 > 0$ such that

$$\mathbb{P} \left[\sum_{i=1}^k (J_{Z^*}(X_i, \lambda_{\bar{X}, Z^*}) - \lambda_{\bar{X}, Z^*} d) \leq \sum_{i=1}^k J_X(X_i, d) + C_2 \log k \right] > 1 - \frac{K_2}{\sqrt{k}} \quad (220)$$

The cdf of the sum of the zero-mean random variables $\lambda_{\bar{X}} (\bar{d}_Z(S_i|X_i = x_i) - \mathbb{E} [\bar{d}_Z(S|X_i)|X_i = x_i])$ is bounded for each $x^k \in \mathcal{T}_k$ as in the proof of (179), leading to the conclusion that there exists $K_1 > 0$ such that for k large enough

$$\mathbb{P} \left[f(S^k, X^k) \leq \sum_{i=1}^k \tilde{f}_{S, X}(S_i, X_i, d) + C_2 \log k \right] > 1 - \frac{K_1 \log k + K_2 + |\mathcal{A}|}{\sqrt{k}} \quad (221)$$

where

$$f(s^k, x^k) \triangleq \sum_{i=1}^k J_{Z^*}(x_i, \lambda_{\bar{X}, Z}) - k\lambda_{\bar{X}, Z}d + \lambda_{\bar{X}} \sum_{i=1}^k (\bar{d}_Z(s_i|x_i) - \mathbb{E}[\bar{d}_Z(S|X_i)|X_i = x_i]) \quad (222)$$

It follows using (215) and (221) that

$$\begin{aligned} \mathbb{P}[g(S^k, X^k) > \log \gamma - \log \beta - k\lambda_0\delta] &\leq \mathbb{P}\left[\sum_{i=1}^k \tilde{J}_{S, X}(S_i, X_i, d) > \log \gamma - \Delta_k\right] \\ &\quad + \frac{K_1 \log k + K_2 + |\mathcal{A}|}{\sqrt{k}} \end{aligned} \quad (223)$$

where $\lambda_0 = \max_{x^k \in \mathcal{T}_k} \lambda_{\bar{X}, Z^*}$ and

$$\Delta_k = \log \beta + k\lambda_0\delta + |\mathcal{A}||\hat{\mathcal{S}}| \log(k+1) + L \log k + C_2 \log k \quad (224)$$

We now weaken the bound in Theorem 4 by choosing

$$\beta = \frac{\sqrt{k}}{b} \quad (225)$$

$$\delta = \frac{\tau}{k} \quad (226)$$

$$\log \gamma = \log M - \log \log_e k + \log 2 \quad (227)$$

where b is that in (207) and $\tau > 0$ is that in (207). Letting

$$\log M = kR(d) + \sqrt{k\tilde{\mathcal{V}}(d)Q^{-1}(\epsilon_k)} + \Delta_k \quad (228)$$

$$\epsilon_k = \epsilon - \frac{K_1 \log k + K_2 + B + \tilde{B} + |\mathcal{A}| + 1}{\sqrt{k}} \quad (229)$$

where \tilde{B} is the Berry-Esséen ratio for $\tilde{J}_{S, X}(S_i, X_i, d)$, and applying (202), (207) and (223), we conclude using Theorem 4 that there exists an (M, d, ϵ') code with M in (228) satisfying

$$\begin{aligned} \epsilon' &\leq \mathbb{P}\left[\sum_{i=1}^k \tilde{J}_{S, X}(S_i, X_i, d) > kR(d) + \sqrt{k\tilde{\mathcal{V}}(d)Q^{-1}(\epsilon_k)}\right] \\ &\quad + \frac{K_1 \log k + K_2 + B + \tilde{B} + |\mathcal{A}| + 1}{\sqrt{k}} \end{aligned} \quad (230)$$

$$\leq \epsilon \quad (231)$$

where (231) is by the Berry-Esséen bound.

REFERENCES

- [1] V. Kostina and S. Verdú, “Nonasymptotic noisy lossy source coding,” in *Proceedings 2013 IEEE Information Theory Workshop*, Seville, Spain, Sep. 2013.
- [2] R. Dobrushin and B. Tsybakov, “Information transmission with additional noise,” *IRE Transactions on Information Theory*, vol. 8, no. 5, pp. 293–304, Sep. 1962.
- [3] T. Berger, *Rate distortion theory*. Prentice-Hall Englewood Cliffs, NJ, 1971.
- [4] H. S. Witsenhausen, “Indirect rate distortion problems,” *IEEE Transactions on Information Theory*, vol. 26, no. 5, pp. 518–521, 1980.
- [5] D. Sakrison, “Source encoding in the presence of random disturbance,” *IEEE Transactions on Information Theory*, vol. 14, no. 1, pp. 165–167, 1968.
- [6] J. Wolf and J. Ziv, “Transmission of noisy information to a noisy receiver with minimum distortion,” *IEEE Transactions on Information Theory*, vol. 16, no. 4, pp. 406–411, 1970.
- [7] E. Ayanoglu, “On optimal quantization of noisy sources,” *IEEE Transactions on Information Theory*, vol. 36, no. 6, pp. 1450–1452, 1990.
- [8] Y. Ephraim and R. M. Gray, “A unified approach for encoding clean and noisy sources by means of waveform and autoregressive model vector quantization,” *IEEE Transactions on Information Theory*, vol. 34, no. 4, pp. 826–834, 1988.
- [9] T. Fischer, J. Gibson, and B. Koo, “Estimation and noisy source coding,” *IEEE Transactions on Acoustics, Speech and Signal Processing*, vol. 38, no. 1, pp. 23–34, 1990.
- [10] T. Courtade and R. Wesel, “Multiterminal source coding with an entropy-based distortion measure,” in *Proceedings 2011 IEEE International Symposium on Information Theory*, Saint-Petersburg, Russia, Aug. 2011, pp. 2040–2044.
- [11] N. Tishby, F. C. Pereira, and W. Bialek, “The information bottleneck method,” in *Proceedings of the 37-th Annual Allerton Conference on Proceedings of the 37-th Annual Allerton Conference on Communication, Control and Computing*, Allerton, IL, Oct. 2000, pp. 368–377.
- [12] V. Kostina and S. Verdú, “Fixed-length lossy compression in the finite blocklength regime,” *IEEE Transactions on Information Theory*, vol. 58, no. 6, pp. 3309–3338, June 2012.
- [13] A. Ingber and Y. Kochman, “The dispersion of lossy source coding,” in *Data Compression Conference (DCC)*, Snowbird, UT, Mar. 2011, pp. 53–62.
- [14] I. Csiszár, “On an extremum problem of information theory,” *Studia Scientiarum Mathematicarum Hungarica*, vol. 9, no. 1, pp. 57–71, Jan. 1974.
- [15] V. Kostina, “Lossy data compression: nonasymptotic fundamental limits,” Ph.D. dissertation, Princeton University, Sep. 2013.
- [16] S. Verdú, “Non-asymptotic achievability bounds in multiuser information theory,” in *50th Annual Allerton Conference on Communication, Control, and Computing*, Allerton, IL, Oct. 2012, pp. 1–8.
- [17] Y. Polyanskiy, “6.441: Information theory lecture notes,” *M.I.T.*, 2012.
- [18] W. Feller, *An Introduction to Probability Theory and its Applications*, 2nd ed. John Wiley & Sons, 1971, vol. II.
- [19] A. Ingber, I. Leibowitz, R. Zamir, and M. Feder, “Distortion lower bounds for finite dimensional joint source-channel coding,” in *Proceedings 2008 IEEE International Symposium on Information Theory*, Toronto, ON, Canada, July 2008, pp. 1183–1187.
- [20] B. Yu and T. Speed, “A rate of convergence result for a universal d-semifaithful code,” *IEEE Transactions on Information Theory*, vol. 39, no. 3, pp. 813–820, 1993.

- [21] I. Csiszár and J. Körner, *Information Theory: Coding Theorems for Discrete Memoryless Systems*, 2nd ed. Cambridge Univ Press, 2011.